

Aggregation of estimators

Olena Verbytska

1 Introduction

An aggregated estimator is a convex or linear combination of estimators. This choice of an estimator allows to preserve all the good properties of each estimator. Such model is less sensitive to small perturbations of the data Y .

All the results are taken from the book [1], Sections 2 and 4. The main result of this paper is Theorem 3.3 about oracle risk bound for the Gibbs mixing estimator which will be compared to oracle risk bound for model selection. Then we consider the numerical approximation of aggregated estimator by Metropolis-Hastings algorithm. And finally we will simulate Gaussian white noise data and apply the theory to recreate the unknown signal.

2 Definitions

Let us consider a set of estimators $\{\hat{f}_m, m \in \mathcal{M}\}$. Instead of selecting only one estimator $\hat{f}_{\hat{m}}$, we will consider a convex combination \hat{f} of the $\{\hat{f}_m\}_m$

$$\hat{f} = \sum_{m \in \mathcal{M}} w_m \hat{f}_m, \quad \text{with} \quad w_m \geq 0 \quad \text{and} \quad \sum_{m \in \mathcal{M}} w_m = 1.$$

Now we consider a collection of models $\{S_m, m \in \mathcal{M}\}$, where $\{S_m\}_m$ are linear subspaces of \mathbb{R}^n . The estimator $\hat{f}_m = \text{Proj}_{S_m} Y$ maximizes the likelihood under the constraint that it belongs to S_m . Then the unbiased estimator of the risk $r_m = \mathbb{E}[\|f^* - \hat{f}_m\|^2]$ is

$$\hat{r}_m = \|Y - \hat{f}_m\|^2 + 2d_m\sigma^2 - n\sigma^2, \quad \text{with } d_m = \dim(S_m).$$

We associate to $\beta > 0$ and a probability distribution π on \mathcal{M} the Gibbs mixing \hat{f} of the collection of estimators $\{\hat{f}_m\}_m$

$$\hat{f} = \sum_{m \in \mathcal{M}} w_m \hat{f}_m, \quad \text{with } w_m = \frac{\pi_m e^{-\beta \hat{r}_m / \sigma^2}}{\mathcal{Z}}, \quad \text{where } \mathcal{Z} = \sum_{m \in \mathcal{M}} \pi_m e^{-\beta \hat{r}_m / \sigma^2}. \quad (1)$$

The Gibbs distribution w corresponds to the distribution *minimizing* over the set of probability q on \mathcal{M} the functional

$$\mathcal{G}(q) = \sum_{m \in \mathcal{M}} q_m \hat{r}_m + \frac{\sigma^2}{\beta} \mathcal{K}(q, \pi), \quad (2)$$

where $\mathcal{K}(q, \pi) = \sum_{m \in \mathcal{M}} q_m \log(q_m / \pi_m) \geq 0$ is the Kullback-Leibler divergence between the probabilities q and π .

3 Theoretical results

Proposition 3.1 (Stein's formula). *Let Y be a n -dimensional Gaussian vector with mean μ and covariance matrix $\sigma^2 I_n$. For any function $F : \mathbb{R}^n \rightarrow \mathbb{R}^n$, $F(x) = (F_1(x), \dots, F_n(x))$ continuous, with piecewise continuous partial derivatives, fulfilling for all $i \in \{1, \dots, n\}$*

(a) $\lim_{|y_i| \rightarrow \infty} F_i(y_1, \dots, y_n) e^{-(y_i - \mu_i)^2 / 2\sigma^2} = 0$ for all $(y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n) \in \mathbb{R}^{n-1}$,

(b) $\mathbb{E} [|\partial_i F_i(Y)|] < \infty$,

we have

$$\mathbb{E} [\|F(Y) - \mu\|^2] = \mathbb{E} [\|F(Y) - Y\|^2 - n\sigma^2 + 2\sigma^2 \operatorname{div}(F)(Y)],$$

where $\operatorname{div}(F) = \sum_i \partial_i F_i$.

Proof.

$$\begin{aligned} \mathbb{E} [\|F(Y) - \mu\|^2] &= \mathbb{E} [\|F(Y) - Y\|^2] + \underbrace{\mathbb{E} [\|Y - \mu\|^2]}_{(2)} + 2 \underbrace{(\mathbb{E} [\langle F(Y), Y - \mu \rangle])}_{(3)} - \underbrace{\mathbb{E} [\langle Y, Y - \mu \rangle]}_{(1)} \\ &= \mathbb{E} [\|F(Y) - Y\|^2] + n\sigma^2 + 2\sigma^2 \mathbb{E} [\operatorname{div}(F)(Y)] - 2n\sigma^2. \end{aligned}$$

Since Y is a Gaussian vector with mean μ and covariance matrix $\sigma^2 I_n$, we have

$$(1) \mathbb{E} [\langle Y, Y - \mu \rangle] = \mathbb{E} [\sum_i (Y_i^2 - Y_i \mu_i)] = \sum_i (\mathbb{E} [Y_i^2] - \mathbb{E} [Y_i] \mu_i) = \sum_i \operatorname{Var} Y_i = n\sigma^2,$$

$$(2) \mathbb{E} [\|Y - \mu\|^2] = \sum_i (\mathbb{E} [Y_i^2] - 2\mathbb{E} [Y_i] \mu_i + \mu_i^2) = \sum_i (\mathbb{E} [Y_i^2] - (\mathbb{E} [Y_i])^2) = \sum_i \operatorname{Var} Y_i = n\sigma^2.$$

An integration by parts, (a) and (b) give

$$\sigma^2 \int_{y_i \in \mathbb{R}} \partial_i F_i(y_1, \dots, y_n) e^{-(y_i - \mu_i)^2 / 2\sigma^2} dy_i = \int_{y_i \in \mathbb{R}} (y_i - \mu_i) F_i(y_1, \dots, y_n) e^{-(y_i - \mu_i)^2 / 2\sigma^2} dy_i,$$

for all $(y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n) \in \mathbb{R}^{n-1}$. Integrating with respect to $y_1, \dots, y_{i-1}, y_{i+1}, \dots, y_n$, it follows that $\mathbb{E} [(Y_i - \mu_i) F_i(Y)] = \sigma^2 \mathbb{E} [\partial_i F_i(Y)]$ and finally we get

$$(3) \mathbb{E} [\langle F(Y), Y - \mu \rangle] = \sigma^2 \mathbb{E} [\operatorname{div}(F)(Y)].$$

□

Now we will compare the oracle risk bounds for model selection and for the Gibbs mixing estimator. In the first case we choose \hat{m} such that for $K > 1$

$$\hat{m} \in \operatorname{argmin}_{m \in \mathcal{M}} \left\{ \|Y - \hat{f}_m\|^2 + \sigma^2 \operatorname{pen}(m) \right\}, \text{ with } \operatorname{pen}(m) = K \left(\sqrt{d_m} + \sqrt{2 \log(1/\pi_m)} \right)^2. \quad (3)$$

Theorem 3.2 (Oracle risk bound for model selection). *For $\hat{f} = \hat{f}_{\hat{m}}$ there exists a constant $C_K > 1$ depending only on $K > 1$, such that*

$$\mathbb{E} [\|\hat{f} - f^*\|^2] \leq C_K \min_{m \in \mathcal{M}} \left\{ \mathbb{E} [\|\hat{f}_m - f^*\|^2] + \sigma^2 \log \left(\frac{1}{\pi_m} \right) + \sigma^2 \right\}. \quad (4)$$

Theorem 3.3 (Oracle risk bound for the Gibbs mixing estimator). *For the Gibbs mixing \hat{f} and $\beta \leq 1/4$ we have*

$$\mathbb{E} [\|\hat{f} - f^*\|^2] \leq \min_{m \in \mathcal{M}} \left\{ \mathbb{E} [\|\hat{f}_m - f^*\|^2] + \frac{\sigma^2}{\beta} \log \left(\frac{1}{\pi_m} \right) \right\}. \quad (5)$$

The risk bounds are very similar. The main difference is that the constant $C_K > 1$ in (4) and $C_K = 1$ in (5). Thus, the Gibbs mixing estimator has better oracle risk bound.

Proof of Theorem 3.3. We apply the Stein's formula to $F(Y) = \hat{f}$ and use the linearity of the divergence operator:

$$\hat{r} = \|Y - \hat{f}\|^2 - n\sigma^2 + 2\sigma^2 \sum_{m \in \mathcal{M}} \operatorname{div} \left(w_m \hat{f}_m \right) \text{ fulfills } \mathbb{E}[\hat{r}] = \mathbb{E} \left[\|\hat{f} - f^*\|^2 \right].$$

After a bounding of \hat{r} , the result will follow from the fact that w minimizes (2).

1. Upper bound of \hat{r} . We first expand

$$\begin{aligned} \|Y - \hat{f}\|^2 &= \left\langle \sum_{m \in \mathcal{M}} w_m (Y - \hat{f}_m), Y - \hat{f}_m + \hat{f}_m - \hat{f} \right\rangle \\ &= \sum_{m \in \mathcal{M}} w_m \langle Y - \hat{f}_m, Y - \hat{f}_m \rangle + \sum_{m \in \mathcal{M}} w_m \langle Y - \hat{f}_m, \hat{f}_m - \hat{f} \rangle \\ &= \sum_{m \in \mathcal{M}} w_m \|Y - \hat{f}_m\|^2 + \underbrace{\sum_{m \in \mathcal{M}} w_m \langle Y - \hat{f}, \hat{f}_m - \hat{f} \rangle}_{=0} - \sum_{m \in \mathcal{M}} w_m \|\hat{f} - \hat{f}_m\|^2, \end{aligned}$$

and notice that the divergence of $w_m \hat{f}_m$ is given by

$$\operatorname{div} \left(w_m \hat{f}_m \right) = w_m \operatorname{div} \left(\hat{f}_m \right) + \left\langle \hat{f}_m, \nabla w_m \right\rangle.$$

Plugging these two formulas in the definition of \hat{r} we obtain

$$\hat{r} = \sum_{m \in \mathcal{M}} w_m \left(\|Y - \hat{f}_m\|^2 - n\sigma^2 + 2\sigma^2 \operatorname{div} \left(\hat{f}_m \right) \right) + \sum_{m \in \mathcal{M}} w_m \left(2\sigma^2 \left\langle \hat{f}_m, \frac{\nabla w_m}{w_m} \right\rangle - \|\hat{f} - \hat{f}_m\|^2 \right).$$

For a linear function $G(Y) = AY$, we have $\partial_i G_i(Y) = A_{ii}$, so $\operatorname{div}(G)(Y) = \operatorname{Tr}(A)$. Since \hat{f}_m is defined as an orthogonal projection $\operatorname{Proj}_{S_m}(Y)$, we have $\operatorname{div}(\hat{f}_m) = \operatorname{Tr}(\operatorname{Proj}_{S_m}) = d_m$ and then

$$\|Y - \hat{f}_m\|^2 - n\sigma^2 + 2\sigma^2 \operatorname{div} \left(\hat{f}_m \right) = \|Y - \hat{f}_m\|^2 - n\sigma^2 + 2d_m\sigma^2 = \hat{r}_m.$$

In addition, we have

$$\begin{aligned} \frac{\nabla w_m}{w_m} &= -2 \frac{\beta}{\sigma^2} (Y - \hat{f}_m) - \frac{\nabla \mathcal{L}}{\mathcal{L}} = -2 \frac{\beta}{\sigma^2} \left((Y - \hat{f}_m) - \sum_{m' \in \mathcal{M}} w_{m'} (Y - \hat{f}_{m'}) \right) \\ &= -2 \frac{\beta}{\sigma^2} (\hat{f} - \hat{f}_m), \end{aligned}$$

$$\begin{aligned}
\hat{r} &= \sum_{m \in \mathcal{M}} w_m \hat{r}_m - 4\beta \sum_{m \in \mathcal{M}} w_m \langle \hat{f}_m, \hat{f} - \hat{f}_m \rangle - \sum_{m \in \mathcal{M}} w_m \|\hat{f} - \hat{f}_m\|^2 \\
&= \sum_{m \in \mathcal{M}} w_m \hat{r}_m - 4\beta \left(\sum_{m \in \mathcal{M}} w_m \langle \hat{f}_m - \hat{f}, \hat{f} - \hat{f}_m \rangle + \underbrace{\sum_{m \in \mathcal{M}} w_m \langle \hat{f}, \hat{f} - \hat{f}_m \rangle}_{=0} \right) - \sum_{m \in \mathcal{M}} w_m \|\hat{f} - \hat{f}_m\|^2 \\
&= \sum_{m \in \mathcal{M}} w_m \hat{r}_m + (4\beta - 1) \sum_{m \in \mathcal{M}} w_m \|\hat{f} - \hat{f}_m\|^2 \\
&\leq \sum_{m \in \mathcal{M}} w_m \hat{r}_m, \text{ when } \beta \leq 1/4.
\end{aligned}$$

2. Optimality condition. Since $\mathcal{K}(w, \pi) \geq 0$ and w minimizes \mathcal{G} over the set of probability on \mathcal{M} , we have

$$\hat{r} \leq \mathcal{G}(w) \leq \mathcal{G}(q) = \sum_{m \in \mathcal{M}} q_m \hat{r}_m + \frac{\sigma^2}{\beta} \mathcal{K}(q, \pi), \quad \text{for any probability } q \text{ on } \mathcal{M}.$$

Taking the expectation, we obtain

$$\mathbb{E} \left[\|\hat{f} - f^*\|^2 \right] = \mathbb{E}[\hat{r}] \leq \mathbb{E}[\mathcal{G}(q)] = \sum_{m \in \mathcal{M}} q_m r_m + \frac{\sigma^2}{\beta} \mathcal{K}(q, \pi)$$

with $r_m = \mathbb{E}[\hat{r}_m] = \mathbb{E} \left[\|\hat{f}_m - f^*\|^2 \right]$. The right-hand side of the above bound is minimum for $q_m = \pi_m e^{-\beta r_m / \sigma^2} / \mathcal{Z}'$ with $\mathcal{Z}' = \sum_m \pi_m e^{-\beta r_m / \sigma^2}$. Since $\mathcal{Z}' \geq \pi_m e^{-\beta r_m / \sigma^2}$ for all $m \in \mathcal{M}$ we observe that

$$\sum_{m \in \mathcal{M}} q_m r_m + \frac{\sigma^2}{\beta} \mathcal{K}(q, \pi) = -\frac{\sigma^2}{\beta} \log(\mathcal{Z}') \leq \min_{m \in \mathcal{M}} \left\{ r_m - \frac{\sigma^2}{\beta} \log(\pi_m) \right\}$$

□

4 Numerical approximation by Metropolis-Hastings

We can compute an approximation of aggregated estimator (1) by using Markov Chain Monte-Carlo (MCMC) technics, which is more efficient than computing the estimator directly, especially when the cardinality of \mathcal{M} is large.

Let $F : \mathcal{M} \rightarrow \mathbb{R}$ be a real function on \mathcal{M} and w be an arbitrary probability distribution on \mathcal{M} fulfilling $w_m > 0$ for all $m \in \mathcal{M}$. The Metropolis-Hastings algorithm is a classical tool to approximate

$$\mathbb{E}_w[F] := \sum_{m \in \mathcal{M}} w_m F(m).$$

The Metropolis-Hastings algorithm allows to generate an ergodic (i.e. aperiodic, irreducible, and positive recurrent) Markov chain $(M_t)_{t \in \mathbb{N}}$ in \mathcal{M} with stationary distribution w and to approximate $\mathbb{E}_w[F]$ by the average $T^{-1} \sum_{t=1}^T F(M_t)$ with T large. This approximation relies on the fact that for any ergodic Markov chain $(M_t)_{t \in \mathbb{N}}$ in \mathcal{M} with stationary distribution w , we have

$$\frac{1}{T} \sum_{t=1}^T F(M_t) \xrightarrow[\text{a.s.}]{T \rightarrow \infty} \mathbb{E}_w[F]. \quad (6)$$

An aperiodic and irreducible Markov chain which is reversible according to the probability w is ergodic. The Metropolis-Hastings algorithm proposes a way to build a Markov chain $(M_t)_{t \in \mathbb{N}}$ the numerical generation of which is computationally efficient.

Let $\Gamma(m, m')$ be the transition probability of an aperiodic irreducible random walk on \mathcal{M} and define

$$\begin{aligned} Q(m, m') &= \Gamma(m, m') \wedge \frac{w_{m'} \Gamma(m', m)}{w_m} \quad \text{for all } m \neq m', \\ Q(m, m) &= 1 - \sum_{m' : m' \neq m} Q(m, m') \quad \text{for all } m \in \mathcal{M}. \end{aligned}$$

Then we can see that $\text{Markov}(w, Q)$ is reversible. Therefore the Markov chain $(M_t)_{t \in \mathbb{N}}$ with transition Q fulfills for any initial condition

$$\frac{1}{T} \sum_{t=1}^T F(M_t) \xrightarrow[\text{a.s.}]{T \rightarrow \infty} \sum_{m \in \mathcal{M}} w_m F(m), \text{ where } F(m) = \hat{f}_m.$$

From a practical point of view, the Markov chain $(M_t)_t$ can be generated as follows. We start from the transition $\Gamma(m, m')$ which is the transition probability of an aperiodic irreducible random walk on \mathcal{M} and then implement the following algorithm.

Algorithm 1: Metropolis-Hastings

Data: a transition probability Γ , an arbitrary $M_1 \in \mathcal{M}$, a burn-in time T_0 ,
 $\{U_t\}_{t=T_0+1}^T \sim \text{Unif}([0, 1]^{T-T_0-1})$.

Result: $\frac{1}{T-T_0} \sum_{t=T_0+1}^T F(M_t)$.

for $t = 1, \dots, T$ **do**

• From the current state M_t generate M'_{t+1} according to the distribution $\Gamma(M_t, \cdot)$;

• $p_{t+1} \leftarrow 1 \wedge \frac{w_{M'_{t+1}} \Gamma(M'_{t+1}, M_t)}{w_{M_t} \Gamma(M_t, M'_{t+1})}$;

if $p_{t+1} < U_t$ (*i.e. with probability p_{t+1} do ...*) **then**

$M_{t+1} \leftarrow M'_{t+1}$.

end

else

$M_{t+1} \leftarrow M_t$.

end

end

The following example shows how to implement the Metropolis-Hastings algorithm for evaluating (1) in the coordinate sparse setting.

Example 4.1. We consider the collection of estimators $\{\hat{f}_m : m \in \mathcal{M}\}$ and the probability π for the coordinate sparse setting. We write $\mathcal{V}(m)$ for the set of all the subsets m' that can be obtained from m by adding or removing one integer $j \in \{1, \dots, p\}$ to m . If we take the uniform distribution on $\mathcal{V}(m)$ as proposal distribution $\Gamma(m, \cdot)$, then $\Gamma(m, m') = \Gamma(m', m)$ for any $m' \in \mathcal{V}(m)$. As a consequence, we only need to compute $w_{M'_{t+1}}/w_{M'_t}$ at the second step of the Metropolis Hastings algorithm. We have seen that

$$\frac{w_{m'}}{w_m} = \frac{\pi_{m'} e^{-\beta \hat{r}_{m'}/\sigma^2}}{\pi_m e^{-\beta \hat{r}_m/\sigma^2}}$$

so we need to compute $\pi_{m'}/\pi_m$ efficiently when $m' \in \mathcal{V}(m)$. Let us consider the choice $\pi_m \propto e^{-|m|}/C_p^{|m|}$ proposed for the coordinate-sparse setting. We have the simple formulas

$$\frac{\pi_{m'}}{\pi_m} = \frac{|m| + 1}{e(p - |m|)} \text{ for } m' = m \cup \{j\} \quad \text{and} \quad \frac{\pi_{m'}}{\pi_m} = \frac{e(p - |m| + 1)}{|m|} \text{ for } m' = m \setminus \{j\}.$$

5 Simulations [2]

We will consider the following Gaussian white noise model:

$$Y = f + \varepsilon, \text{ where } f(x) = x^5 - 8x^3 + 10x + 6, \quad \varepsilon \sim N(0, \sigma^2 I_n).$$

We suppose, that σ^2 is known and equals to 4.

```
1 import numpy as np
2 import matplotlib.pyplot as plt
3 from numpy.random import default_rng
4
5 rng = default_rng()
6 n = 100
7 f = lambda x: x**5 - 8*x**3 + 10*x + 6
8 seq = np.linspace(-2.5, 2.5, n)
9 sig2 = 4
10 eps = rng.normal(0, sig2, size=n)
11 Y = f(seq) + eps
```

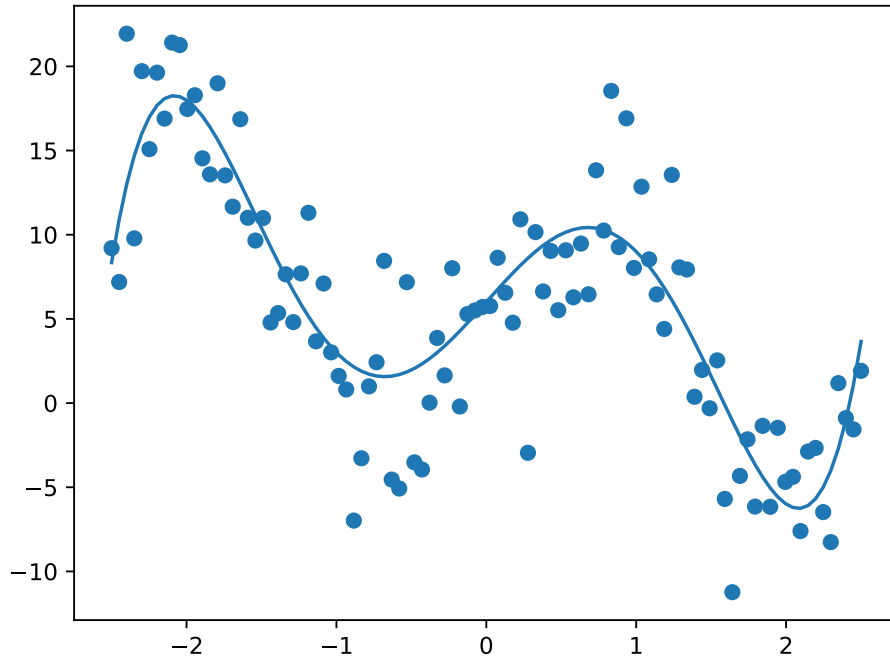


Figure 1: The simulated data Y (blue dots) and the unknown signal f (blue line).

We want to estimate this signal by expanding the observations on the Fourier basis.

$$\left\{ \begin{array}{l} \varphi_1(x) = \frac{1}{\sqrt{P}}, \end{array} \right. \quad (7)$$

$$\left\{ \begin{array}{l} \varphi_{2k}(x) = \sqrt{\frac{2}{P}} \cos\left(\frac{2\pi kx}{P}\right), \end{array} \right. \quad (8)$$

$$\left\{ \begin{array}{l} \varphi_{2k+1}(x) = \sqrt{\frac{2}{P}} \sin\left(\frac{2\pi kx}{P}\right), \text{ where } P \text{ is a period.} \end{array} \right. \quad (9)$$

```

1 # Fourier basis matrix (n_fct is odd)
2 def create_fourier_matrix(seq, n_fct):
3     n = len(seq)
4     basis_matrix = np.zeros((n, n_fct))
5     basis_matrix[:, 0] = np.ones(n) / np.sqrt(100)
6
7     for k in range(int((n_fct-1)/2)):
8         basis_matrix[:, 2*k+1] = np.sin(2 * np.pi * (k+1) * seq / 5) / np.sqrt(50)
9         basis_matrix[:, 2*k+2] = np.cos(2 * np.pi * (k+1) * seq / 5) / np.sqrt(50)
10
11     return(basis_matrix)

```

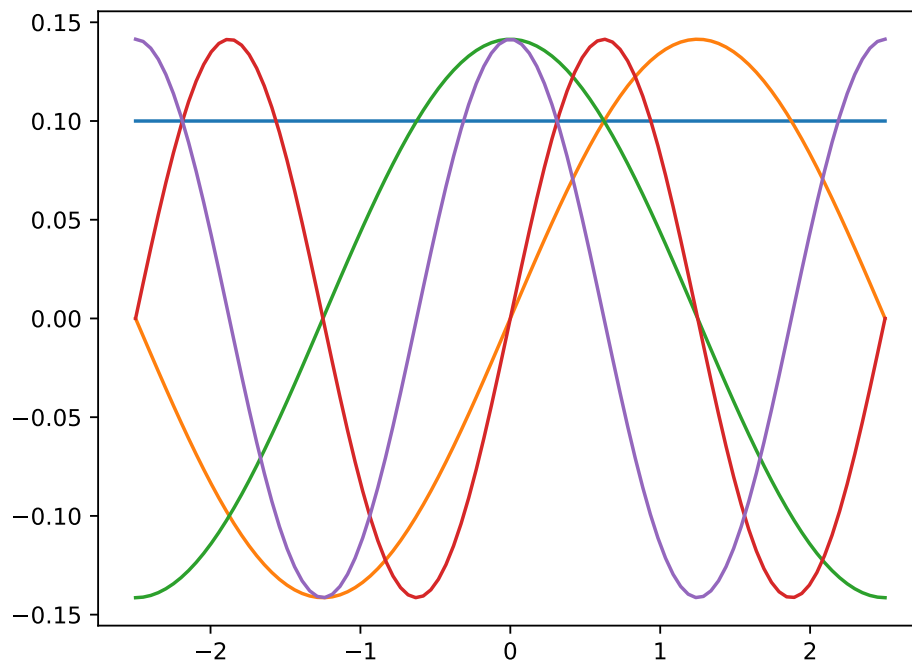


Figure 2: The first 5 functions of Fourier basis

The matrix $X \in \mathbb{R}^{n \times p}$ can be considered as orthogonal, because $X^T X \approx I_p$.

```

1 X = create_fourier_matrix(seq, 9)
2 np.round(X.T @ X, 2)
3
4 '''array([[ 1.   ,  0.   , -0.01,  0.   ,  0.01,  0.   , -0.01,  0.   ,  0.01],
5          [ 0.   ,  0.99, -0.   ,  0.   , -0.   ,  0.   , -0.   ,  0.   ,  0.   ],
6          [-0.01, -0.   ,  1.01,  0.   , -0.02,  0.   ,  0.02, -0.   , -0.02],
7          [ 0.   ,  0.   ,  0.   ,  0.99,  0.   , -0.   ,  0.   , -0.   ,  0.   ],
8          [ 0.01, -0.   , -0.02,  0.   ,  1.01, -0.   , -0.02,  0.   ,  0.02],
9          [ 0.   ,  0.   ,  0.   , -0.   , -0.   ,  0.99,  0.   ,  0.   ,  0.   ],
10         [-0.01, -0.   ,  0.02,  0.   , -0.02,  0.   ,  1.01, -0.   , -0.02],
11         [ 0.   ,  0.   , -0.   , -0.   ,  0.   ,  0.   , -0.   ,  0.99,  0.   ],
12         [ 0.01,  0.   , -0.02,  0.   ,  0.02,  0.   , -0.02,  0.   ,  1.01]])'''

```


Now we will calculate the coefficients Z of Y in the Fourier basis: $Z = X^T Y$, $Z \sim N(\theta, \sigma^2 I_p)$. Analysing the positioning of Z with respect to the lines $y = \pm\sqrt{2\sigma^2 \log(p)}$, we can consider the expectation θ as sparse vector.

```

1 p = 101
2 X = create_fourier_matrix(seq, p)
3 Z = X.T @ Y
4 plt.scatter(np.linspace(-2.5, 2.5, len(Z)), Z)
5 plt.axline((0, np.sqrt(2 * sig2 * np.log(p))), (1, np.sqrt(2 * sig2 * np.log(p)
6   )), color='r')
7 plt.axline((0, -np.sqrt(2 * sig2 * np.log(p))), (1, -np.sqrt(2 * sig2 * np.log(
8   p))), color='r')
9 plt.show()

```

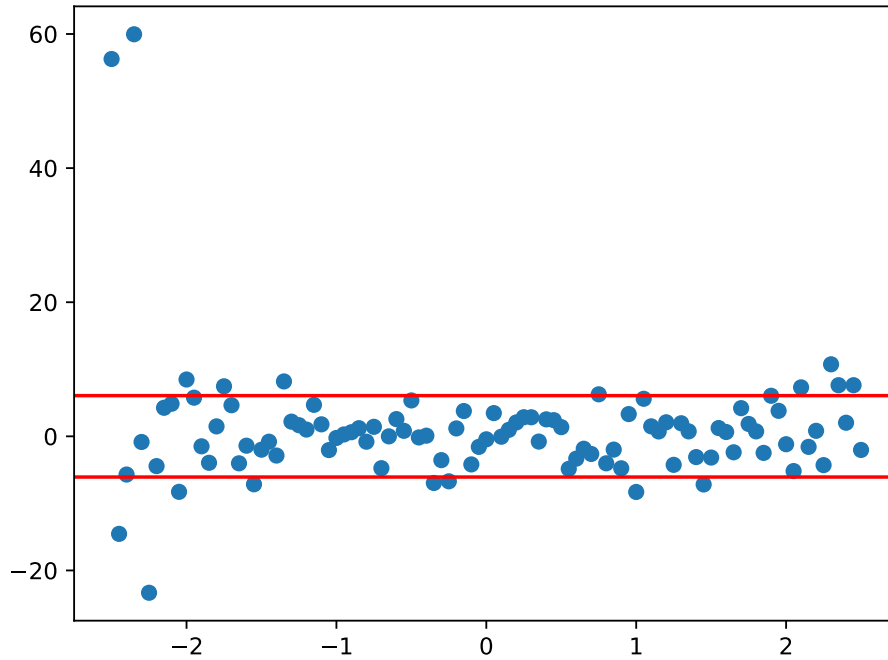


Figure 3: Coefficients Z of Y in the Fourier basis

We implemented the thresholding estimator for θ :

$$\hat{\theta}_j = Z_j \mathbb{1}_{\{|Z_j| \geq \sqrt{2\sigma^2 \log(p)}\}}, \text{ for } j = 1, \dots, p.$$

```

1 n_models = 50
2 error = []
3 for k in range(n_models):
4     X = create_fourier_matrix(seq, (2*k+1))
5     Z = X.T @ Y
6     p = 2*k+1
7     t = np.sqrt(2 * sig2 * np.log(p))
8     theta_hat = Z * (np.abs(Z) >= t)
9     f_hat = X @ theta_hat

```

```

10 error.append(np.sum((f_hat - f(seq))**2))
11
12 k_opt = 2*(error.index(min(error))+1)+1
13 print("m_hat =", k_opt)
14 X = create_fourier_matrix(seq, k_opt)
15 Z = X.T @ Y
16 p = k_opt
17 t = np.sqrt(2 * sig2 * np.log(p))
18 theta_hat = Z * (np.abs(Z) >= t)
19 f_hat = X @ theta_hat
20 plt.scatter(seq, Y)
21 plt.plot(seq, f(seq))
22 plt.plot(seq, f_hat, color='orange')
23 print("Error =", np.round(np.sum((f_hat - f(seq))**2), 2))

```

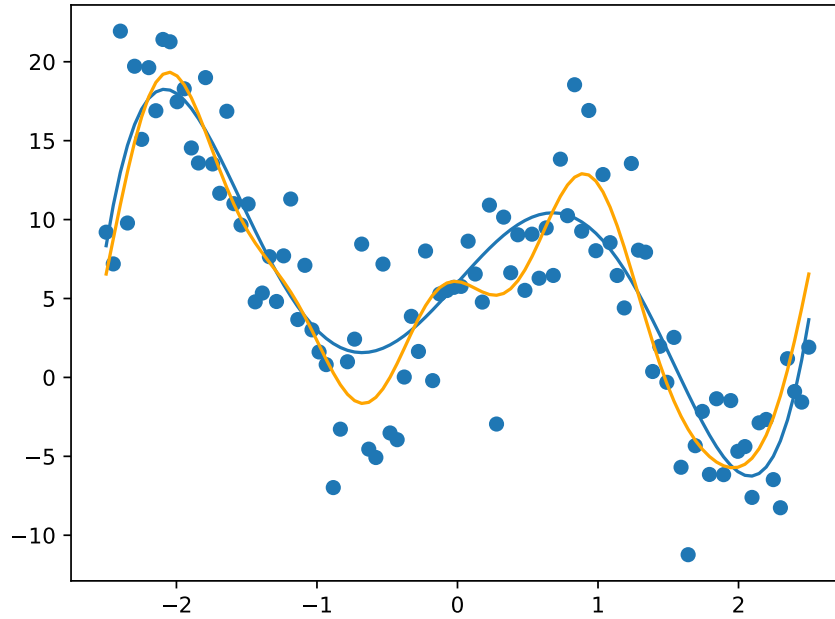


Figure 4: The selection model with the best oracle risk, error = 357.04.

Probably, the choice of another basis could have given better results.

We take all π_m 's equal and select \hat{m} according to the criterion (3).

```

1 n_models = 50
2 K = 1.1
3 error = []
4 for k in range(n_models):
5     X = create_fourier_matrix(seq, (2*k+1))
6     Z = X.T @ Y
7     p = 2*k+1
8     t = np.sqrt(2 * sig2 * np.log(p))
9     theta_hat = Z * (np.abs(Z) >= t)
10    f_hat = X @ theta_hat

```

```

11 m_hat = np.sum((f_hat - Y)**2) + sig2 * K * (1 + np.sqrt(2 * np.log(n_models)
12 ))**2
13 error.append(m_hat)
14 k_opt = 2*(error.index(min(error))+1)+1
15 print("m_hat =", k_opt)
16 X = create_fourier_matrix(seq, k_opt)
17 Z = X.T @ Y
18 p = k_opt
19 t = np.sqrt(2 * 4 * np.log(p))
20 theta_hat = Z * (np.abs(Z) >= t)
21 f_hat = X @ theta_hat
22 print("Error =", np.round(np.sum((f_hat - f(seq))**2), 2))

```

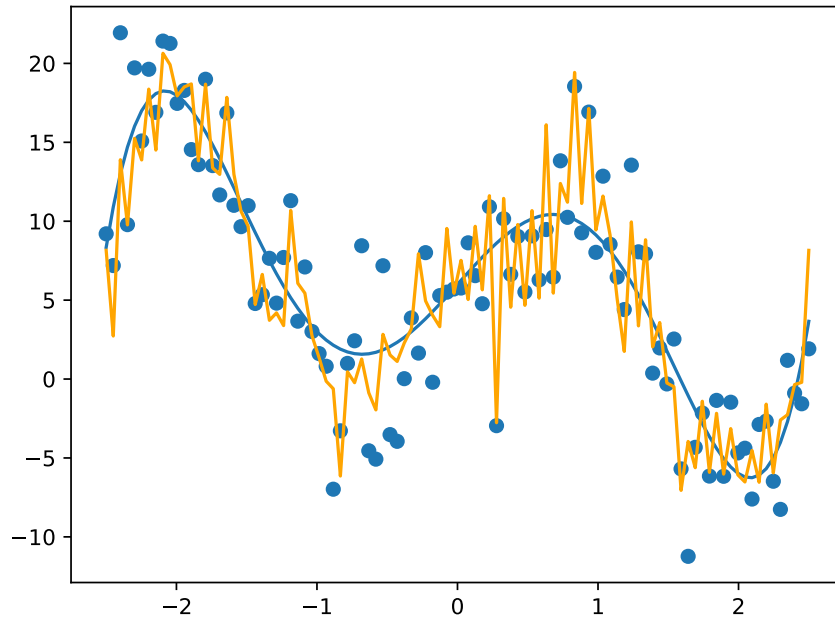


Figure 5: The selection model \hat{m} , error = 1079.9. Overfitting.

Now let us try to approximate F by Metropolis-Hastings algorithm. $\Gamma(M_t, \cdot)$ is uniform on M .

```

1 n_models = 50
2 # A set of estimators f_hat_set and a set of risks r_hat_set
3 f_hat_set = np.array([])
4 r_hat_set = np.array([])
5 for k in range(n_models):
6     X = create_fourier_matrix(seq, (2*k+1))
7     Z = X.T @ Y
8     p = 2*k+1
9     t = np.sqrt(2 * 4 * np.log(p))
10    theta_hat = Z * (np.abs(Z) >= t)
11    f_hat = X @ theta_hat
12    f_hat_set = np.append(f_hat_set, f_hat)
13    dim = 2 * k + 1
14    r_hat = np.sum((Y-f_hat)**2) + 2 * dim * sig2 - n * sig2

```

```

15 r_hat_set = np.append(r_hat_set, r_hat)
16 f_hat_set = f_hat_set.reshape(n_models, n)
17
18 T_0 = 0
19 T = 10000
20 M_t = 4
21 U = rng.uniform(size=T)
22 F = np.zeros(n)
23 beta = 1/4
24
25 for t in range(T-T_0):
26     M_prime = rng.choice(n_models, 1) # uniform choice of M'_{t+1} from M
27     if M_prime - M_t == 1:
28         ratio_pi = (M_t + 1) / (np.exp(1) * (n_models - M_t))
29     elif M_t - M_prime == 1:
30         ratio_pi = (np.exp(1) * (n_models - M_t + 1)) / M_t
31     else:
32         ratio_pi = 1 # for easier calculations
33     ratio_w = ratio_pi * np.exp(beta / sig2 * (r_hat_set[M_prime] - r_hat_set[M_t]
34         ]))
35     p_t = np.min([1, ratio_w])
36     if p_t < U[t]:
37         M_t = M_prime
38     else:
39         M_t = M_t
40     F = F + f_hat_set[M_t, :] / (T-T_0)

```

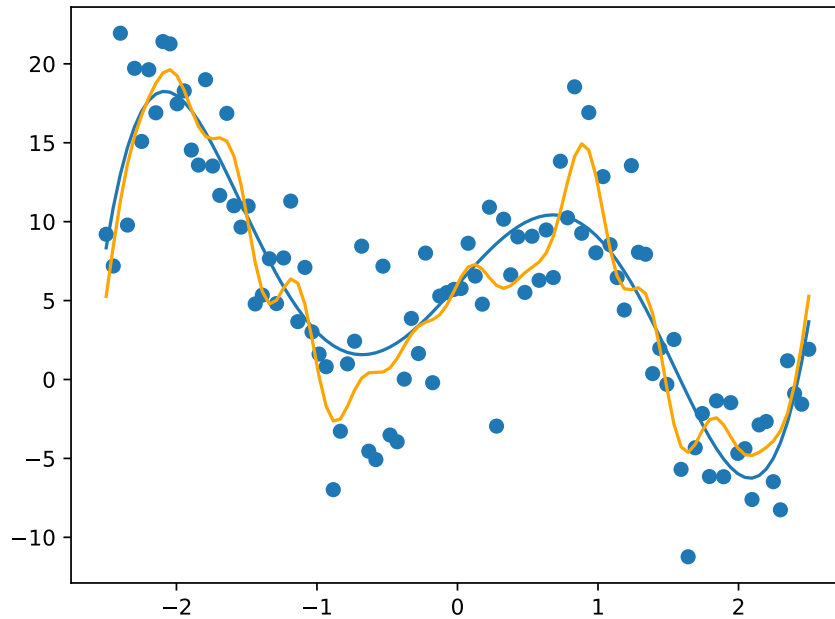


Figure 6: The model with aggregated estimator \hat{f} , error = 437.85.

As we can see, on this type of data aggregated estimator indeed performs better than the selected one.

References

- [1] Christophe Giraud. *Introduction to High-Dimensional Statistics*. Universite Paris Saclay, 2021.
- [2] *Google colab: Aggregation of estimators.ipynb*. URL: <https://colab.research.google.com/drive/1iWqtSil1MifM6GlW-rRaVHvZop8aEPdZ?usp=sharing>.